





HANDS ON SCRAPING DATA DENGAN KASUS ANALISIS SENTIMEN PADA KEBIJAKAN PEMERINTAH

Wahyu Calvin Frans Mariel – BPS Muhammad Ghozy Al Haqqoni – BPS

Rabu, 28 Mei 2025



Web Scraping?



- ✓ Teknik mengumpulkan data yang jumlahnya bisa sedikit hingga sangat besar.
- ✓ Berisi berbagai informasi yang disediakan oleh suatu website.
- ✓ Biasanya dilakukan setelah proses pengumpulan URL dari web crawling.



Web Crawling?



- ✓ Suatu proses **menjelajahi halaman** suatu website yang biasanya memprogram untuk indexing seluruh web.
- ✓ Bisa digunakan untuk meningkatkan kualitas SEO (search engine optimization) suatu web. Seperti Googlebot, Skyscanner, webarchive.

Objective

- 1. Netiquette: Semua website boleh di-scraping?
- 2. Mengenal lebih dalam: Website
- 3. Tools untuk web-scraping
- **4.Hands-on**: Real world cases



Netiquette: Web Scraping Ethics

Penting!

Konsiderasi memilih *website* yang dapat dilakukan *web scraping*

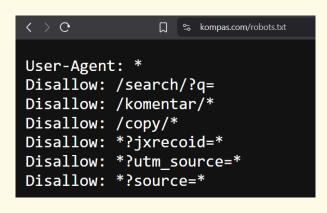
- 1. Term and Conditions Websites
- 2. Robots.txt (easier way)
- 3. Cloudflare atau anti-bot software lainnya



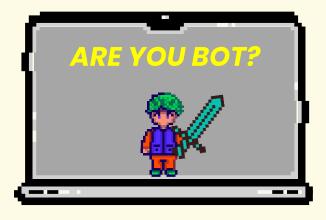
Netiquette: Robots.txt

File robots.txt biasanya digunakan oleh pemilik *website* untuk mengendalikan akses dan memberikan instruksi pada bot mesin pencari tentang *rules* apa yang diperbolehkan dan dilarang karena alasan privasi atau keamanan.

Bisa cek dari https://(websiteyangdituju)/robots.txt



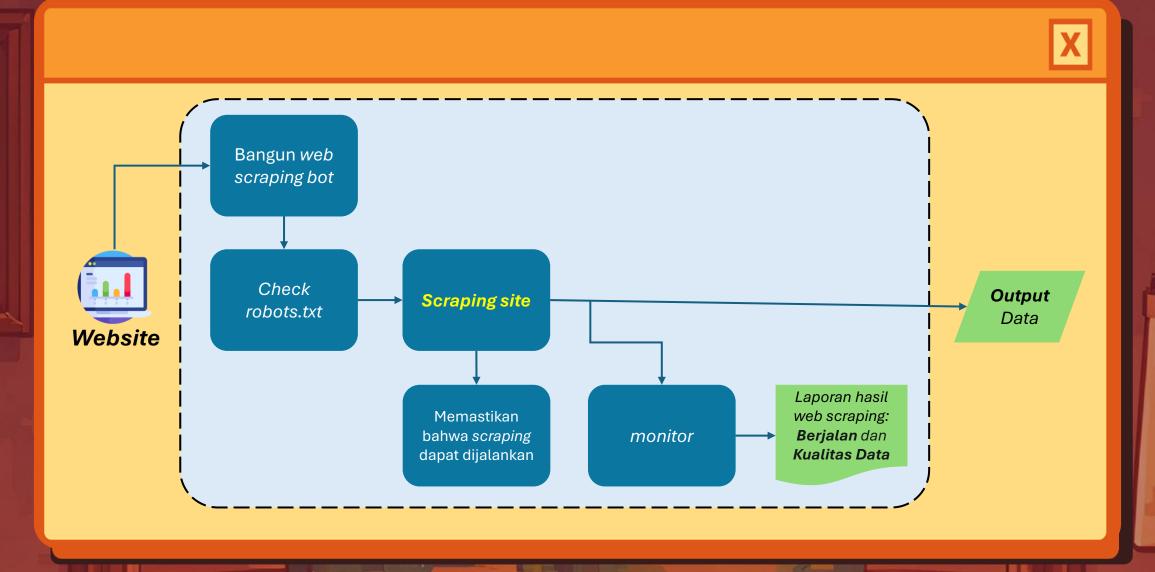




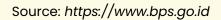


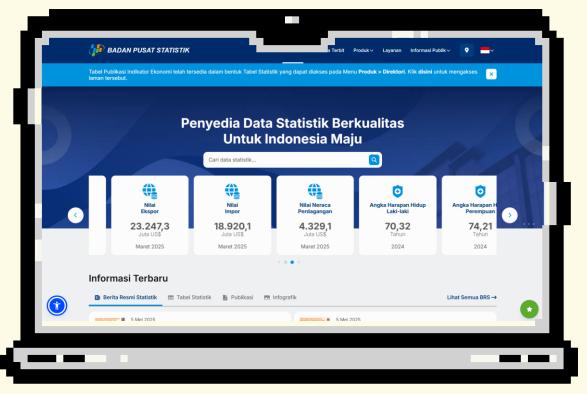


Guidance to Web Scraping



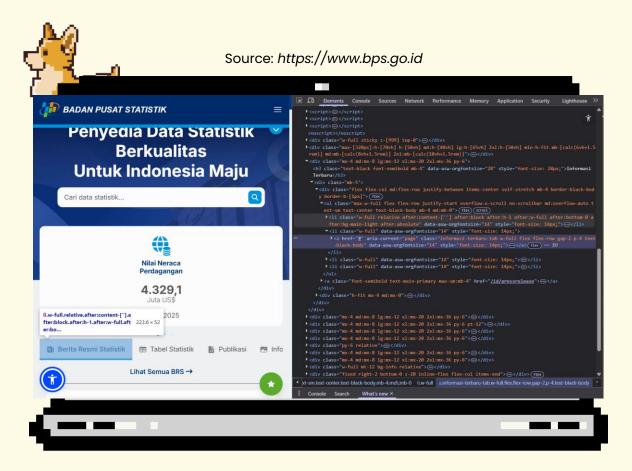
Website







Website

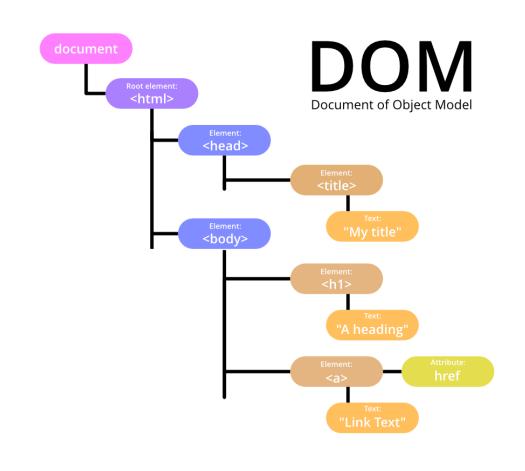


Belakang layar

Dibangun dengan *programming language* dan di-*load* melalui HTML



HyperText Markup Language



Struktur dari suatu *website*

- ✓ Membangun dan menghubungkan suatu web pages satu dan lainnya
- ✓ ,Menampilkan berbagai informasi: head and body



Source: https://www.geeksforgeeks.org/what-isdocument-object-in-java-dom/

HTML: Tags

Tag	Description	
<html> </html>	Declares the Web page to be written in HTML	
<head> </head>	Delimits the page's head	
<title> </title>	Defines the title (not displayed on the page)	
<body> </body>	Delimits the page's body	
<h<i>n> </h<i> n>	Delimits a level <i>n</i> heading	
 	Set in boldface	
<i> </i>	Set in italics	
<center> </center>	Center on the page horizontally	
 	Brackets an unordered (bulleted) list	
 	Brackets a numbered list	
 	Brackets an item in an ordered or numbered list	
	Forces a line break here	
<	Starts a paragraph	
<hr/>	Inserts a horizontal rule	
	Displays an image here	
 	Defines a hyperlink	

Simbol untuk menandai suatu elemen di HTML

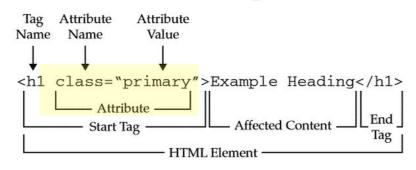
- ✓ Berisi instruksi kepada browser bagaimana suatu objek bisa ditampilkan
- ✓ objek bisa berupa text, video, audio, dan gambar



Source: https://medium.com/@nwachukwufavourc/understandinghtml-tags-and-elements-795402e3eb47

Tags Attributes

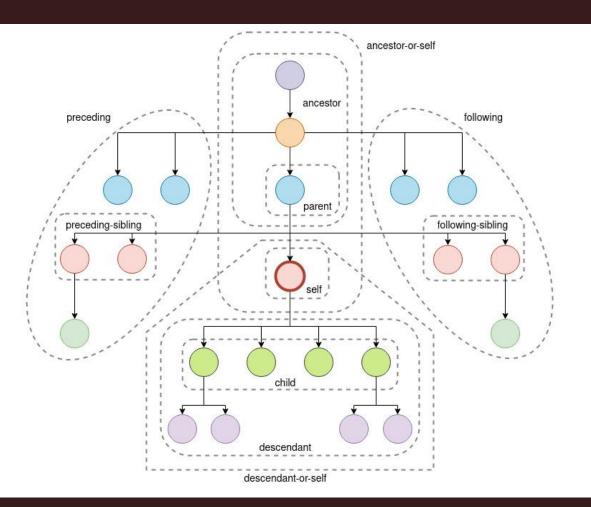
HTML Tags



Source: https://medium.com/@nwachukwufavourc/understandinghtml-tags-and-elements-795402e3eb47

Komponen untuk mengidentifikasi sebuah tag

Tags Axes



Source: https://jrebecchi.github.io/xpathhelper/xpath-axes.html#parent-axis/

Identifikasi *tag*s berdasarkan hubungan antar elemen

- ✓ Menavigasi hierarki atas/bawah tags
- ✓ Parent-to-sibling, Parent-to-child, dsb.

```
<html>
    <head>
        <title>Website iprakom</title>
        </head>
        <body>
            <h1>Selamat datang</h1>
            Jadwal Sharing Knowledge:

                  25 Mei 2025
                  27 Mei 2025 (Today!)
                  </body>
                 </html>
```

Contoh:

li: child to ul

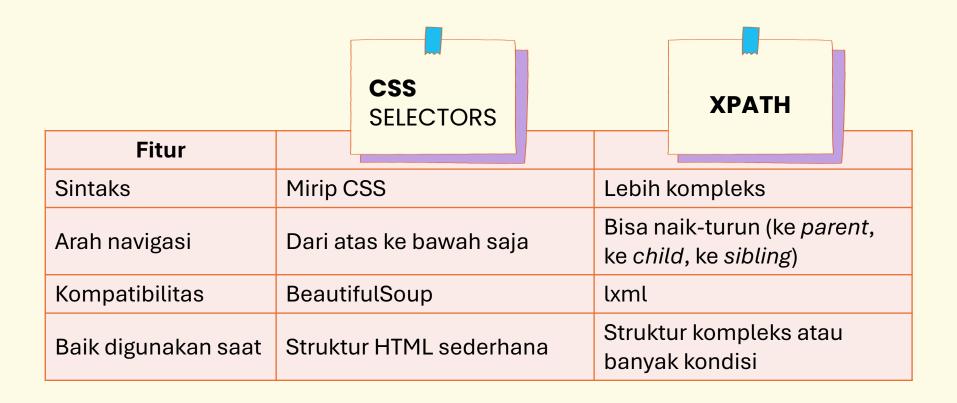
p: sibling to h1 and ul

The Way to Select Tag(s)

CSS SELECTORS

XPATH

The Way to Select Tag(s)



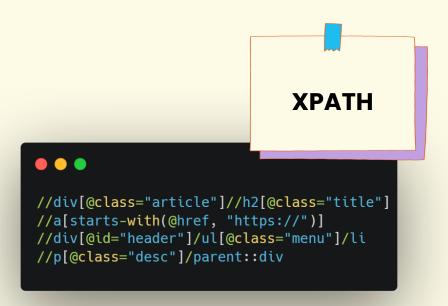
The Way to Select Tag(s)

PROMPT:

- 1. ambil <h2 class="title"> dalam <div class="article">
- 2. semua <a> dengan href diawali https://
- 3. elemen di dalam dalam <div id="header">
- 4. elemen <div class="container"> yang punya anak

css
SELECTORS

div.article h2.title
a[href^="https://"]
div#header > ul.menu > li



Web Scraping Tools

- **✓ MANUALLY**
- ✓ BANTUAN PIHAK KETIGA
- **✓ GOOGLE SHEETS**
- **✓ JAVASCRIPT**
- ✓ R
- **✓ PYTHON**
- ✓ AND SO ON...



Python Libraries













Type of Web Scraping

HTML Extraction

Mengumpulkan data dengan mengekstraksi informasi yang ada dalam struktur HTML web page

API Extraction

Mengumpulkan data melalui API (*Application Programming Interface*) yang disediakan suatu *website.* Biasanya *output* dalam JSON

Bisa dilakukan pengecekan pada tab network di developer tools browser:

Windows: Ctrl + Shift + I

MacOS: Command + Option + I ($\Re + \neg + I$)



LET'S RUN A QUICK SIMULATION



Kalau mendengar **kebijakan pemerintah**, yang terpikirkan oleh teman-teman, kita bisa **memperoleh informasi** terkait kebijakan itu **dari mana** saja?



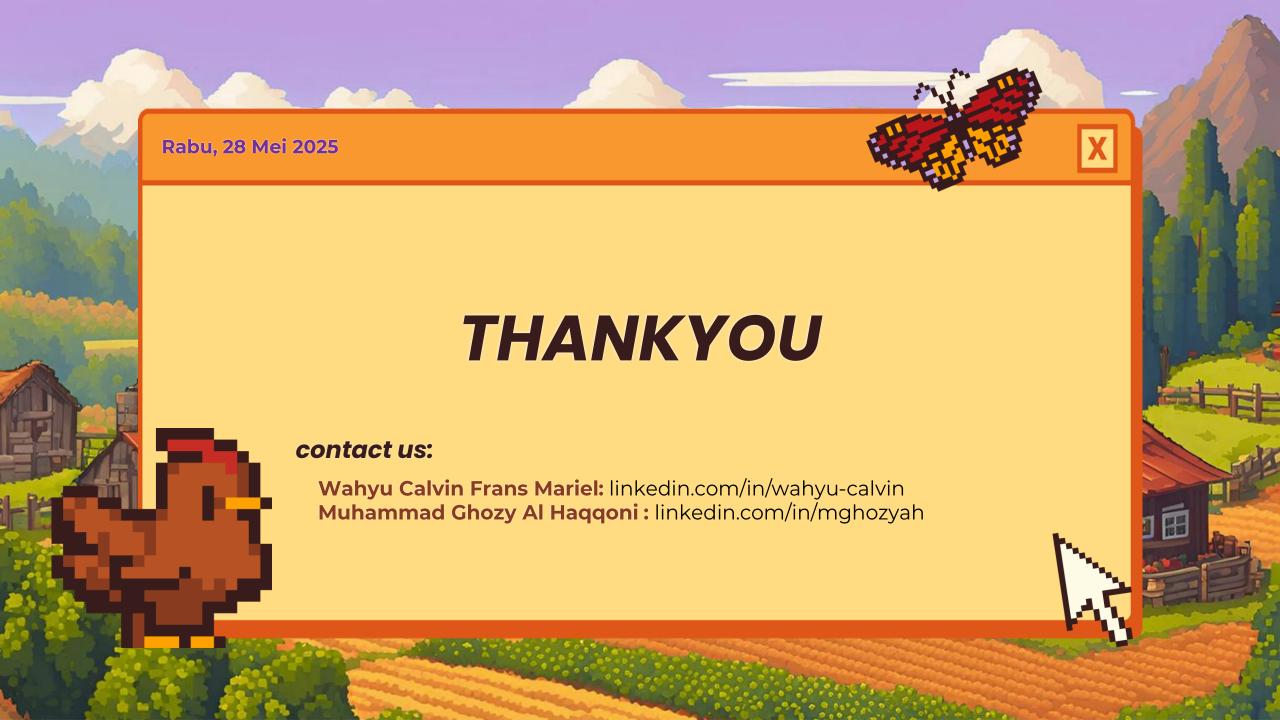
SOSIAL MEDIA BERITA YOUTUBE WHATSAPP CHANNEL DENGER LANGSUNG











Lampiran: The Way to Select Tag(s)

CSS SELECTORS

Tujuan	Contoh CSS Selector	Penjelasan	JLLC
Pilih elemen berdasarkan tag	div	Semua <div></div>	
Berdasarkan kelas	.container	Elemen dengan class="container"	
Berdasarkan id	#header	Elemen dengan id="header"	
Tag + kelas	div.article	<div> dengan class="article"</div>	
Tag + id	div#main	<div> dengan id="main"</div>	
Atribut tertentu	a[href="https://"]	<a> dengan href persis "https://"	
Atribut yang dimulai dengan	a[href^="https"]	href mulai dengan "https"	
Child langsung	div > p	yang langsung anak <div></div>	
Semua descendant	div p	Semua di dalam <div></div>	
Sibling langsung	h2 + p	langsung setelah <h2></h2>	
Semua sibling	h2 ~ p	Semua setelah <h2> pada level yar</h2>	ng sama

Lampiran: The Way to Select Tag(s)

XPATH

Tujuan	Contoh XPath	Penjelasan
Semua elemen tag tertentu	//div	Semua <div> di dokumen</div>
Berdasarkan atribut kelas	//div[@class="container"]	<div> dengan class="container"</div>
Berdasarkan atribut id	//div[@id="header"]	<div> dengan id="header"</div>
Tag + atribut tertentu	//a[@href="https://"]	<a> dengan href persis "https://"
Atribut yang dimulai dengan	//a[starts-with(@href, "https")]	href mulai dengan "https"
Child langsung	/div/p	yang langsung anak <div></div>
Semua descendant	//div//p	Semua di dalam <div></div>
Naik ke parent	//p[@class="desc"]/parent::div	Parent <div> dari</div>
Pilih berdasarkan posisi	//ul/li[3]	ketiga dalam
Pilih berdasarkan teks	//a[contains(text(), "Berita")]	<a> yang mengandung kata "Berita"